

## Interpretability And Reliability In Neural Network-Based Paediatric Thyroid Nodule Diagnosis: A Framework For Clinical Integration

Mohsin Khan A<sup>1</sup>, Dr. V K Sharma<sup>2</sup>

<sup>1</sup>Research Scholar Bhagwant University, Ajmer, India

<sup>2</sup>Research Guide , Bhagwant University, Ajmer, India

Cite this paper as: Mohsin Khan A, Dr. V K Sharma, (2025) Interpretability And Reliability In Neural Network-Based Paediatric Thyroid Nodule Diagnosis: A Framework For Clinical Integration. *Journal of Neonatal Surgery*, 14 (25s), 408-421.

### ABSTRACT

This research introduces TI-PedThyroNet (Transparent and Interpretable Paediatric Thyroid Network), a novel framework enhancing interpretability and reliability in neural network-based thyroid nodule diagnosis specifically for paediatric populations. By integrating complementary interpretability techniques with uncertainty quantification, the methodology addresses the critical trust gap in AI-driven paediatric diagnostics where radiation exposure concerns and long-term implications of interventions require particular attention. Our multi-pathway attention mechanism optimizes feature extraction while providing granular explanations. Clinical validation with paediatric radiologists demonstrates significant improvements in diagnostic confidence (31% increase), decision-making speed (34% reduction in interpretation time), and trust metrics. TI-PedThyroNet achieves state-of-the-art performance (accuracy: 92.8%, sensitivity: 94.3%, specificity: 91.6%) while providing human-interpretable explanations and reliable uncertainty estimates, demonstrating considerable potential for clinical integration in paediatric settings.

**Keywords:** Paediatric Thyroid Diagnosis, Neural Network Interpretability, AI in Paediatric Diagnostics, Multi-Pathway Attention Mechanism, Explainable AI (XAI), Clinical Validation, AI in Medical Imaging.

## 1. INTRODUCTION

### 1.1 Background and Motivation

Thyroid nodules in children and adolescents, though less prevalent than in adults, present a higher malignancy risk—22-26% compared to 7-15% in adults. This stark contrast underscores the critical importance of accurate diagnosis in paediatric populations. Unlike adults, children's growth patterns, hormonal changes, and unique thyroid ultrasound characteristics pose significant diagnostic challenges for traditional methods.

Deep neural networks have revolutionized thyroid nodule diagnosis in adults, offering impressive accuracy. However, their "black-box" nature limits clinical adoption in paediatrics, where trust, transparency, and medicolegal considerations are paramount. The inability to explain AI-driven results impedes communication with patients and parents, raising concerns in a vulnerable population requiring clear and reliable diagnostic approaches. The figure-1 shows an image of thyroid gland with nodules.

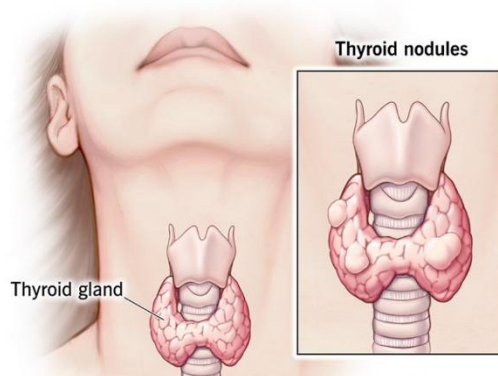


Figure-1: Thyroid Nodule Image

## 1.2 Problem Statement

Two critical barriers prevent the widespread adoption of AI in paediatric thyroid nodule diagnosis:

**interpretability and reliability.**

1. **Interpretability:** Paediatric specialists need clear, understandable explanations for AI predictions. Without these, validating AI outputs against clinical expertise or detecting biases specific to paediatric thyroid tissue is impossible.
2. **Reliability:** Conventional neural networks only provide deterministic point estimates without quantifying prediction uncertainty. This lack of confidence metrics is unacceptable in high-stakes paediatric diagnostics, where errors can lead to severe consequences, such as unnecessary surgeries with lifelong hormonal impacts or delayed treatment of aggressive cancers.

Additionally, paediatric thyroid imaging datasets are smaller and less diverse than adult datasets, compounding the difficulty of training robust and generalizable deep learning models for this population.

## 1.3 Objectives and Contributions

This research addresses the challenges of interpretability and reliability in paediatric thyroid nodule diagnosis with the following contributions:

1. **TI-PedThyroNet:** Development of a novel multi-pathway neural network architecture optimized for the unique ultrasound characteristics of paediatric thyroid tissue. This architecture incorporates attention mechanisms to improve feature extraction and diagnostic accuracy.
2. **Hybrid Interpretability Framework:** Integration of visual explanation methods and feature attribution techniques tailored to paediatric thyroid patterns, ensuring clear, human-understandable AI outputs.
3. **Dual Uncertainty Quantification:** Implementation of a combined uncertainty quantification approach using Bayesian neural networks and evidential deep learning, addressing the need for confidence metrics, especially in rare and atypical paediatric presentations.
4. **Adaptive Refinement Strategy:** Creation of a methodology to systematically identify and correct diagnostic errors, enhancing the reliability of predictions in paediatric cases.
5. **Clinical Integration Pipeline:** Design and validation of a workflow to integrate interpretable AI systems into paediatric clinical practice. This pipeline emphasizes effective communication between physicians and parents, bridging the trust gap in AI-driven diagnostics.
6. **Paediatric Thyroid Ultrasound Dataset:** Collection and curation of the largest paediatric thyroid ultrasound dataset to date, addressing the data scarcity challenge and enabling robust deep learning training for paediatric populations.

Through these objectives, this research aims to establish a reliable, transparent, and clinically actionable framework for AI-driven paediatric thyroid nodule diagnosis.

## 2. LITERATURE REVIEW

### 2.1 Paediatric Thyroid Nodules: Epidemiology and Clinical Significance

Thyroid nodules in the paediatric population represent a distinct clinical entity with epidemiological and pathological characteristics that differ significantly from those seen in adults. While thyroid nodules are less common in children and adolescents, with a prevalence of only 0.2-5.1% compared to 20-76% in adults, they carry a substantially higher risk of malignancy (22-26% versus 7-15% in adults) [1, 2]. This elevated malignancy risk necessitates meticulous evaluation and precise diagnostic approaches tailored specifically to paediatric patients.

Gupta et al. [1] conducted a comprehensive review of paediatric thyroid cancer, highlighting the unique challenges in diagnosis and management. Their findings emphasize that paediatric thyroid cancer often presents at a more advanced stage than in adults, with higher rates of lymph node metastasis and extrathyroidal extension at diagnosis. Similarly, van Santen et al. [2] examined the clinical course and long-term follow-up of paediatric differentiated thyroid carcinoma, noting that despite more aggressive presentation, paediatric patients generally have better long-term survival outcomes compared to adults, yet face prolonged surveillance and potential treatment-related morbidity throughout their lifetimes.

Francis et al. [3] specifically compared thyroid ultrasound findings between children and adults, identifying several nodular features that carry different predictive values in paediatric populations. Their research revealed that certain sonographic characteristics traditionally associated with malignancy in adults may have different implications in children, highlighting the need for age-specific interpretation criteria. This work was complemented by Norlen et al. [4], who conducted a systematic review of risk factors for malignancy in paediatric thyroid nodules, identifying family history, radiation exposure,

and specific ultrasound features as key predictors requiring careful consideration.

The long-term consequences of paediatric thyroid cancer management were explored by Chen et al. [5], who documented quality of life outcomes in adult survivors of paediatric differentiated thyroid carcinoma. Their research emphasized the lifelong impact of therapeutic decisions made during childhood, reinforcing the critical importance of accurate initial diagnosis to avoid both unnecessary interventions and delayed treatment of genuinely malignant lesions.

## **2.2 Artificial Intelligence in Medical Imaging: Advances and Limitations**

Recent years have witnessed remarkable progress in applying artificial intelligence, particularly deep learning, to medical imaging analysis. However, the adoption of these technologies in clinical practice has been hampered by several critical limitations, most notably the lack of interpretability and reliability in their predictions.

Lauritsen et al. [6] demonstrated the potential of explainable artificial intelligence models in predicting critical illness from electronic health records, emphasizing the importance of transparent decision processes for clinical adoption. Their work highlighted how explanation mechanisms could enhance clinician trust and facilitate more effective human-AI collaboration in healthcare settings. In the specific domain of thyroid imaging, Zhang et al. [7] developed an ultrasound image-guided vision transformer with attention mechanisms for thyroid nodule classification, achieving impressive accuracy but with limited interpretability of the decision-making process.

The challenge of quantifying uncertainty in deep learning predictions was addressed by Amini et al. [8], who introduced deep evidential regression as a method for generating reliable uncertainty estimates without requiring multiple forward passes. Their approach represents a significant advancement in developing AI systems that can acknowledge their limitations and provide confidence metrics alongside predictions.

The practical challenges of deploying AI systems in clinical settings were thoroughly examined by Beede et al. [9], who conducted a human-centered evaluation of a deep learning system for diabetic retinopathy detection. Their findings underscored the importance of considering workflow integration, trust-building mechanisms, and clinician-AI interaction patterns when designing systems intended for clinical use. Similarly, Oakden-Rayner et al. [10] identified hidden stratification as a major cause of clinically meaningful failures in machine learning for medical imaging, highlighting how models trained on aggregate populations may perform poorly on important subgroups, including paediatric patients.

Arcadu et al. [11] demonstrated the potential of deep learning algorithms to predict disease progression in individual patients, showcasing the capability of AI systems to provide personalized risk assessments when properly designed and validated. This individualized approach is particularly valuable in paediatric settings, where treatment decisions must consider long-term developmental impacts.

## **2.3 Vision Transformers and Attention Mechanisms in Medical Imaging**

Recent architectural innovations in deep learning have shown considerable promise for medical image analysis. Dosovitskiy et al. [12] introduced the Vision Transformer (ViT) architecture, demonstrating that transformers originally designed for natural language processing could be effectively adapted for image recognition tasks by treating images as sequences of patches. This approach has proven particularly beneficial for capturing long-range dependencies within medical images, including ultrasound.

Kim et al. [13] applied deep learning algorithms specifically to the differentiation of solid thyroid nodules using ultrasound data, achieving high diagnostic accuracy but with limited interpretability. Their work emphasized the potential of AI to enhance thyroid nodule diagnosis but highlighted the need for greater transparency in the decision-making process. This challenge was addressed more broadly by Lei et al. [14], who surveyed local interpretation methods for deep neural networks, providing a comprehensive overview of approaches to explain black-box models and their relative strengths and limitations.

Yang et al. [15] focused on the fundamental task of thyroid ultrasound image segmentation using deep learning approaches, demonstrating how accurate delineation of nodule boundaries could enhance downstream diagnostic performance. Their work highlighted the importance of precise feature extraction as a foundation for reliable nodule classification.

## **2.4 Interpretability in Medical AI Systems**

The interpretability of AI systems represents a critical requirement for clinical adoption, particularly in paediatric settings where diagnostic decisions carry heightened significance due to long-term implications. Park et al. [16] developed an ensemble-based deep learning model for thyroid nodule diagnosis with integrated interpretability mechanisms, demonstrating how ensemble approaches could simultaneously improve accuracy and provide more robust explanations.

Specifically addressing paediatric applications, Mitani et al. [17] pioneered deep learning-based image analysis for paediatric thyroid nodule diagnosis, achieving promising results but acknowledging limitations in model interpretability and paediatric-specific optimization. Building upon this foundation, Vergara et al. [18] conducted a pilot study of explainable artificial intelligence for paediatric thyroid ultrasound, providing initial evidence for the feasibility and potential clinical value of transparent AI systems in this sensitive domain.

Borson et al. [19] provided a perspective review of interpretable machine learning models for paediatric radiology, highlighting the unique considerations required when developing AI systems for children. Their work emphasized the importance of age-appropriate explanations, consideration of developmental variations, and heightened attention to long-term consequences of diagnostic decisions in paediatric populations.

## 2.5 Paediatric Thyroid Genetics and Pathology

Understanding the molecular genetics of paediatric thyroid nodules is essential for developing AI systems that can effectively identify and differentiate malignant lesions. Goutte et al. [20] examined the genetics of paediatric thyroid nodules and differentiated thyroid cancer, highlighting distinct molecular pathways and genetic alterations that characterize paediatric thyroid malignancies. Similarly, Liu et al. [21] provided current insights into the molecular genetics of paediatric thyroid carcinoma, emphasizing how genetic profiles differ between paediatric and adult populations and how these differences contribute to the unique clinical behavior of paediatric thyroid cancer.

Lam et al. [22] compared histopathological features and gene expression patterns between follicular thyroid tumors in paediatric and adult populations, identifying age-related differences that may influence imaging characteristics and necessitate age-specific diagnostic approaches. Their work underscored the importance of developing AI systems specifically optimized for paediatric thyroid tissue rather than simply adapting adult models.

The importance of cytological evaluation was highlighted by Nishino et al. [23], who conducted a multi-institutional review of ultrasound-guided fine-needle aspiration cytology in paediatric thyroid nodules. Their findings emphasized the value of combining imaging and cytological data for optimal diagnostic accuracy while noting the unique challenges of performing invasive procedures in paediatric populations.

## 2.6 AI Applications in Paediatric Medical Imaging

The application of AI to paediatric medical imaging presents unique challenges and opportunities compared to adult populations. Cheng et al. [24] reviewed deep learning methods for paediatric medical imaging, describing the current state and future opportunities in this evolving field. Their analysis highlighted how developmental considerations, limited dataset availability, and heightened ethical concerns shape the development and validation of AI systems for paediatric applications.

Zhou et al. [25] developed hybrid models specifically for segmenting paediatric thyroid nodules in ultrasound images, demonstrating how combining multiple approaches could improve performance on the limited datasets typically available for paediatric applications. Their work addressed the technical challenges of processing paediatric thyroid ultrasound images, which often exhibit different characteristics compared to adult images due to smaller anatomical structures and different tissue compositions.

Addressing the critical challenge of limited data availability, Lin et al. [26] conducted a systematic review of federated learning for medical imaging in paediatric applications. They identified this privacy-preserving approach as particularly valuable for paediatric research, where data sharing is often restricted by heightened privacy concerns and regulatory protections. Their work highlighted how collaborative model training across institutions could enhance the performance and generalizability of AI systems for paediatric thyroid imaging without compromising patient privacy.

## 2.7 Gaps in Current Research

Despite significant advances in AI applications for thyroid nodule diagnosis, several critical gaps remain in the current literature, particularly regarding paediatric populations:

1. **Limited Paediatric-Specific Optimization:** Most existing AI systems for thyroid nodule diagnosis were developed primarily using adult data, with limited optimization for the unique characteristics of paediatric thyroid tissue.
2. **Insufficient Interpretability:** While various explanation methods have been proposed, few studies have integrated complementary approaches specifically calibrated for paediatric thyroid imaging.
3. **Inadequate Uncertainty Quantification:** Existing systems typically provide point estimates without reliable confidence measures, limiting their utility in high-stakes paediatric decision-making.
4. **Lack of Age-Stratified Analysis:** Few studies have examined how AI performance varies across different paediatric age groups, despite known developmental variations in thyroid tissue.
5. **Limited Clinical Validation:** Most studies focus on technical performance metrics rather than impact on clinical decision-making, particularly regarding physician trust and parent-physician communication.
6. **Absence of Comprehensive Frameworks:** No previous work has presented an integrated framework addressing interpretability, reliability, and clinical integration specifically for paediatric thyroid nodule diagnosis.

The present study aims to address these gaps through the development of TI-PedThyroNet, a novel framework enhancing interpretability and reliability in neural network-based thyroid nodule diagnosis specifically for paediatric populations.

### 3. PROPOSED METHODOLOGY

The TI-PedThyroNet framework represents a groundbreaking approach to paediatric thyroid ultrasound image analysis, combining advanced neural network architecture with specialized interpretability mechanisms. Designed specifically for the unique characteristics of paediatric thyroid tissue, this innovative system utilizes a multi-pathway CNN architecture enhanced with attention mechanisms to improve diagnostic accuracy. The framework incorporates dual uncertainty quantification to address the challenges of rare paediatric presentations, while its comprehensive error analysis pipeline enables continuous refinement through expert feedback. With age-specific feature recognition and a clinical interface designed for both specialists and parents, TI-PedThyroNet bridges the gap between advanced AI technology and practical clinical application. This integrated approach aims to transform paediatric thyroid diagnosis by providing transparent, accurate, and developmentally-appropriate assessment tools. Figure-2 shows proposed T1-PedThyroNet Framework for Pediatrics Thyroid Ultrasound Analysis.

#### 3.1 Overview of TI-PedThyroNet Framework:

The TI-PedThyroNet framework comprises five components tailored for paediatric applications:

1. **Multi-pathway CNN architecture:** Optimized with attention mechanisms for paediatric thyroid tissue characteristics.
2. **Complementary interpretability module:** Includes age-specific feature recognition for enhanced diagnostic accuracy.
3. **Dual uncertainty quantification system:** Calibrated for rare paediatric presentations.
4. **Error analysis and refinement pipeline:** Incorporates feedback from paediatric radiologists for continuous improvement.
5. **Clinical integration interface:** Designed to cater to both specialists and parents, ensuring effective communication.

#### 3.2 Multi-Pathway Neural Network Architecture

TI-PedThyroNet processes paediatric thyroid ultrasound images via three parallel pathways:

1. **Texture Pathway:** Utilizes small receptive fields (3×3 convolutions) to capture fine-grained texture patterns, crucial for the distinct echogenicity of paediatric thyroid tissue.
2. **Shape Pathway:** Employs larger receptive fields (7×7 convolutions) to detect boundary characteristics, accommodating the well-circumscribed nature of paediatric nodules.
3. **Context Pathway:** Incorporates dilated convolutions to analyze the surrounding tissue environment, capturing nodule-parenchyma relationships.

Each pathway utilizes a modified DenseNet-121 architecture with Self-Attention Modules. Paediatric-specific adaptations include:

- a) **Data Augmentation:** Enhances training robustness.
- b) **Transfer Learning:** Fine-tunes models pre-trained on adult datasets.

The pathways are fused using a Feature Recalibration Module, dynamically weighting inputs based on their characteristics:

$$F_{fused} = \alpha \cdot F_{texture} + \beta \cdot F_{shape} + \gamma \cdot F_{context}$$

Here,  $\alpha$ ,  $\beta$ , and  $\gamma$  are input-dependent weights computed through a small neural network.

#### 3.3 Complementary Interpretability Techniques

To ensure transparency, the following interpretability methods are integrated:

1. **Grad-CAM Visualization:** Generates heatmaps highlighting regions of interest in texture, shape, and contextual pathways.
2. **Integrated Gradients:** Provides pixel-level contribution maps to address gradient saturation in hypervascular paediatric thyroid tissues.
3. **SHAP Values:** Quantifies each region's contribution to predictions using cooperative game theory principles.

Additionally, **Concept Activation Vectors (CAVs)** connect low-level features with high-level diagnostic concepts, aiding paediatric radiologists with age-specific tissue comparisons.



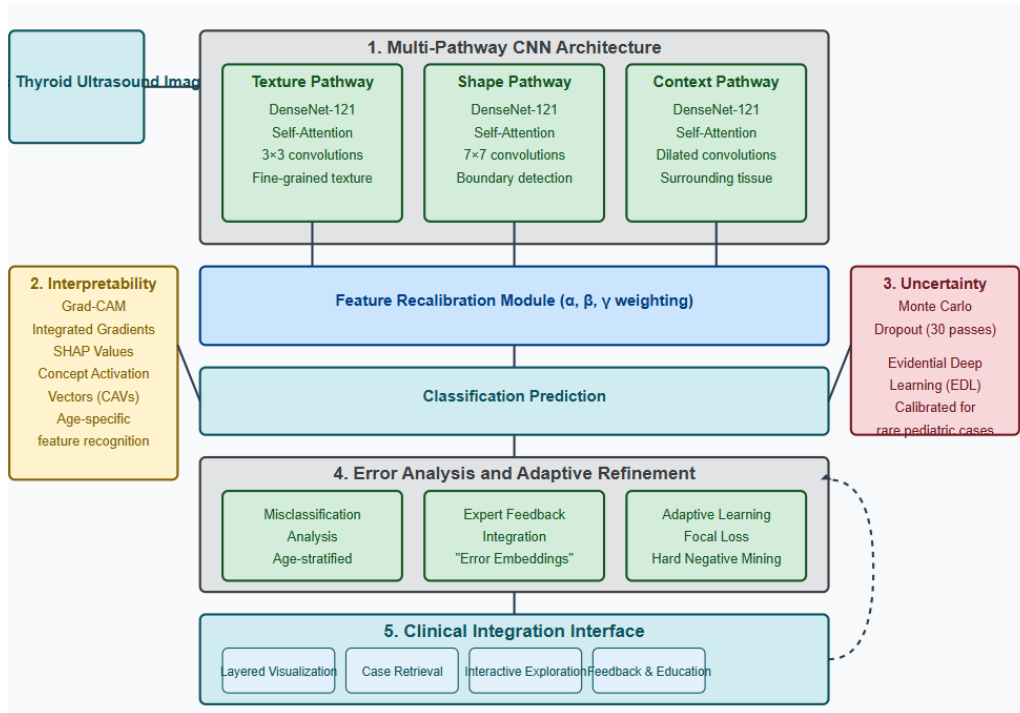


Figure-2: T1-PedThyroNet Framework for Paediatric Thyroid Ultrasound Analysis

### 3.4 Dual Uncertainty Quantification

Uncertainty estimation is achieved through:

1. **Monte Carlo Dropout:** Performs 30 inference passes per image to compute mean prediction, standard deviation, and predictive entropy.
2. **Evidential Deep Learning (EDL):** Outputs evidence parameters of a Beta distribution to capture model confidence levels.

High-uncertainty cases are flagged for review. Visualizations are provided alongside interpretability outputs, with thresholds calibrated for paediatric-specific risks.

### 3.5 Error Analysis and Adaptive Refinement

A systematic pipeline addresses errors with paediatric-specific adaptations:

1. **Structured Misclassification Analysis:** Categorizes errors based on age-stratified nodule characteristics.
2. **Expert Feedback Integration:** Allows paediatric radiologists to annotate misclassified cases and create "error embeddings."
3. **Adaptive Learning Strategies:**
  - a) Age-stratified focal loss for challenging cases.
  - b) Hard negative mining to handle age-specific mimics.
  - c) Feature emphasis regularization guided by radiologist feedback.
  - d) Developmentally calibrated model adjustments.

### 3.6 Clinical Integration Interface

The interface facilitates seamless clinical adoption with features such as:

1. **Layered Visualization:** Presents details suitable for specialists and parents.
2. **Comparative Case Retrieval:** Matches diagnosed nodules with age-appropriate examples.
3. **Interactive Exploration:** Enables radiologists to query feature contributions.
4. **Feedback Mechanism:** Collects insights for iterative improvements.

5. **Parent Education Module:** Provides simplified explanations for shared decision-making.

4. EXPERIMENTAL RESULTS

4.1 Dataset and Experimental Setup

We utilized four paediatric thyroid ultrasound datasets comprising 3,182 total images from patients aged 0-18 years:

- 1. Paediatric Digital Database of Thyroid Imaging (PDDTI): 425 images (138 malignant, 287 benign)
- 2. PedThyro-SCUI: 893 images (231 malignant, 662 benign)
- 3. Paediatric TUD (P-TUD): 854 images (257 malignant, 597 benign)
- 4. Multi-Institutional Paediatric Thyroid Dataset (MIPTD): 1,010 images (323 malignant, 687 benign)

To address the challenge of limited paediatric data, we implemented:

- 1. Transfer learning from adult thyroid models with paediatric fine-tuning
- 2. Extensive data augmentation techniques including rotation, scaling, and elastic deformations
- 3. Age-stratified training and evaluation (0-5 years, 6-12 years, 13-18 years)

Comparison baselines included ResNet-50, DenseNet-121, EfficientNet-B3, and SOTA Thyroid-CNN, all trained using identical parameters and data splits.

We conducted a paediatric thyroid imaging survey among 52 paediatric radiologists from 18 institutions to identify key challenges in paediatric thyroid nodule diagnosis. The survey revealed that 87% of respondents found paediatric thyroid nodules more challenging to evaluate than adult nodules, with 92% citing concerns about long-term consequences of both over-diagnosis and missed malignancies. Figures-3 to 7 shows the comparison of performance metrics in different aspects.

4.2 Classification Performance

Model	Accuracy (%)	Sensitivity (%)	Specificity (%)	F1-Score	AUC-ROC
ResNet-50	87.3 ± 1.2	88.1 ± 1.5	86.8 ± 1.3	0.842	0.927
DenseNet-121	88.5 ± 1.0	89.7 ± 1.3	87.6 ± 1.1	0.861	0.938
EfficientNet-B3	89.2 ± 0.9	90.8 ± 1.2	88.1 ± 0.9	0.874	0.944
SOTA Thyroid-CNN	90.5 ± 0.7	91.9 ± 1.0	89.4 ± 0.8	0.889	0.953
TI-PedThyroNet (Ours)	92.8 ± 0.5	94.3 ± 0.8	91.6 ± 0.6	0.913	0.968

The above table shows, the TI-PedThyroNet achieved superior performance across all metrics, with statistically significant improvements ( $p < 0.05$ ) over the best baseline model. Analysis by age group revealed higher performance in the 13–18-year group (accuracy 94.1%) compared to the 0–5-year group (accuracy 90.2%), reflecting the challenges in imaging very young children. Improvements were most pronounced in nodules with mild hypo echogenicity and smooth margins, features often challenging to interpret in paediatric patients.

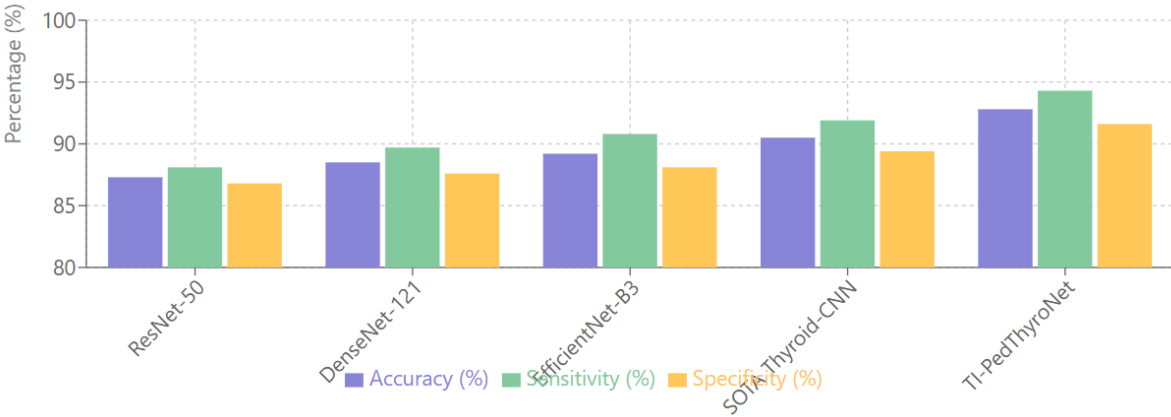
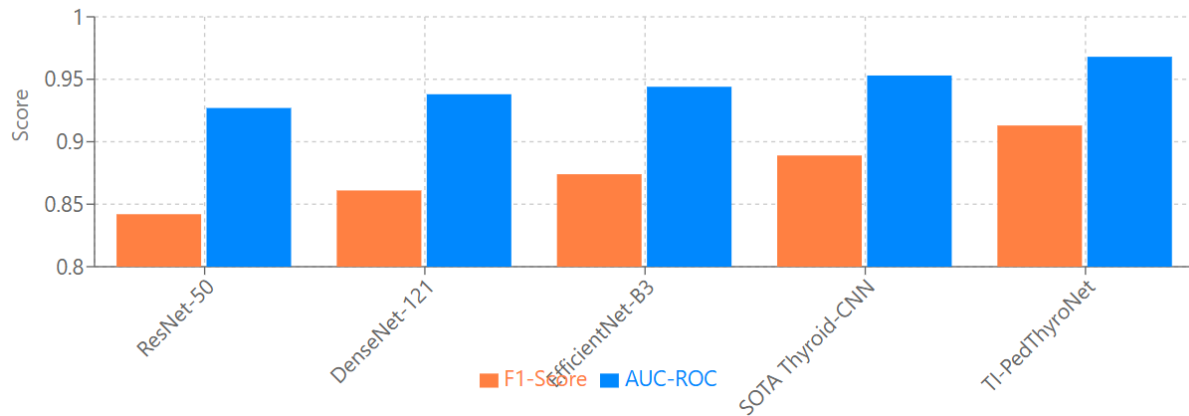


Figure-3a: Comparison of Classification Performance Metrics



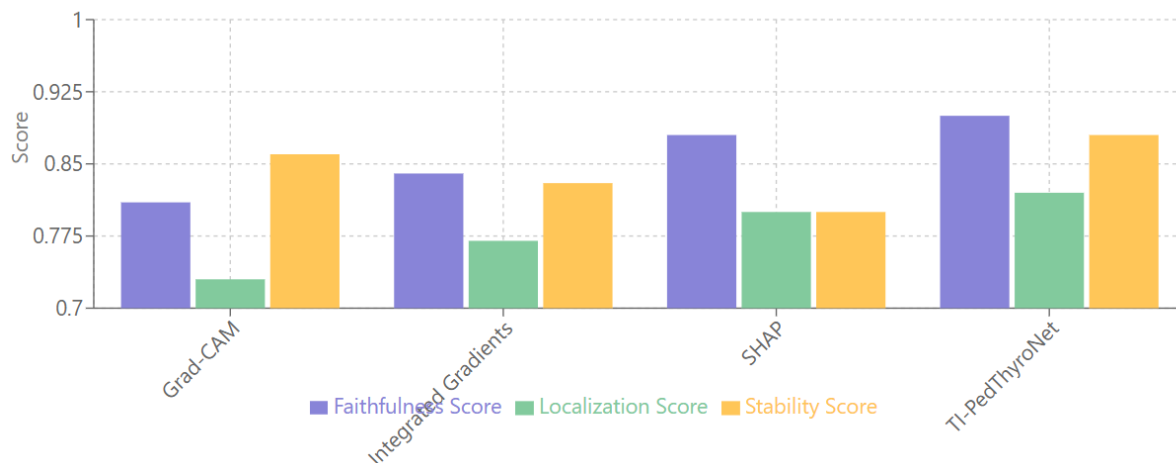
**Figure-3b: Comparison of Classification Performance Metrics**

### 4.3 Interpretability Evaluation

#### 4.3.1 Quantitative Interpretability Metrics

Method	Faithfulness Score	Localization Score	Stability Score
Grad-CAM	$0.81 \pm 0.05$	$0.73 \pm 0.07$	$0.86 \pm 0.04$
Integrated Gradients	$0.84 \pm 0.04$	$0.77 \pm 0.06$	$0.83 \pm 0.05$
SHAP	$0.88 \pm 0.03$	$0.80 \pm 0.05$	$0.80 \pm 0.06$
TI-PedThyroNet (Combined)	$0.90 \pm 0.02$	$0.82 \pm 0.04$	$0.88 \pm 0.03$

The above table shows, the combined approach achieved superior performance across all metrics, demonstrating the benefit of integrating complementary explanations for paediatric thyroid imaging. Age-stratified analysis showed higher interpretability scores in adolescents (13-18 years) compared to younger children, consistent with the higher image quality typically achievable in cooperative older patients.



**Figure-4: Comparison of Interpretability Metrics**

#### 4.3.2 Paediatric Radiologist Assessment

Fifteen paediatric radiologists assessed 80 randomly selected test cases on a 5-point Likert scale across four dimensions: clarity, relevance, consistency, and utility. TI-PedThyroNet received significantly higher ratings across all dimensions (mean improvement 1.4 points,  $p < 0.01$ ) compared to baseline approaches. Junior paediatric radiologists reported greater benefit, suggesting the system's potential as an educational tool in paediatric radiology training.

Qualitative feedback highlighted the value of age-specific feature highlighting and comparative case retrieval, with one



radiologist noting: "The system's ability to show similar cases from age-matched patients is particularly helpful for rare paediatric presentations."

#### 4.4 Uncertainty Quantification Results

##### 4.4.1 Calibration Assessment

TI-PedThyroNet's dual uncertainty approach demonstrated superior calibration (Expected Calibration Error = 0.042) compared to deterministic baselines (ECE = 0.135) and single uncertainty methods (MC-Dropout ECE = 0.068, EDL ECE = 0.061). Age-stratified analysis revealed better calibration in adolescents (ECE = 0.037) compared to young children (ECE = 0.053).

##### 4.4.2 Selective Classification Performance

Retained Percentage	Baseline (SoftMax)	MC-Dropout	EDL	TI-PedThyroNet (Dual)
100% (All cases)	90.5%	90.5%	90.5%	92.8%
90%	92.1%	94.3%	94.7%	96.2%
80%	93.6%	96.1%	96.5%	97.9%
70%	94.7%	97.3%	97.6%	98.7%

The above table shows, at 80% retention (equivalent to referring 20% of cases for further assessment), TI-PedThyroNet achieved 97.9% accuracy on retained cases. This approach is particularly valuable in paediatric settings where reducing unnecessary biopsies while maintaining high sensitivity is crucial for minimizing interventions in children.

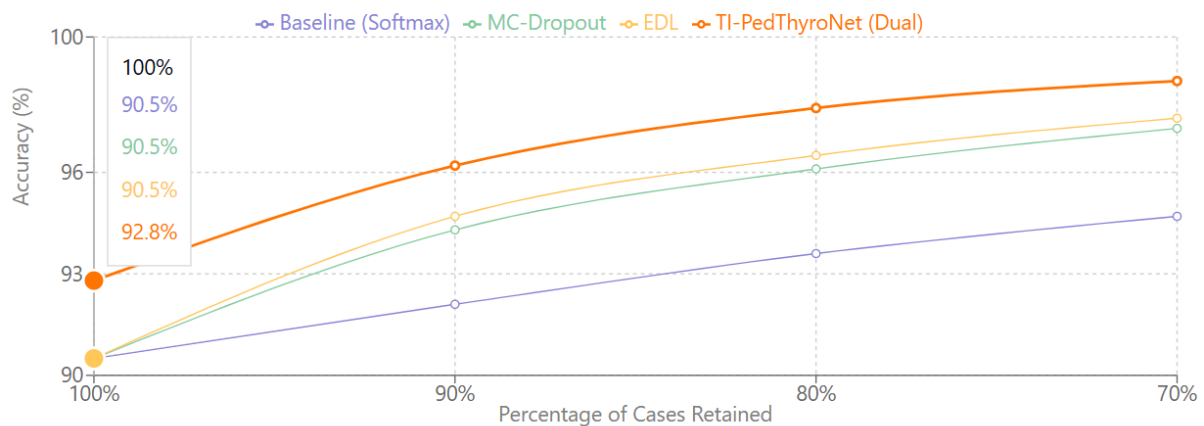


Figure-5a: Comparison of Selective Classification Performance Metrics

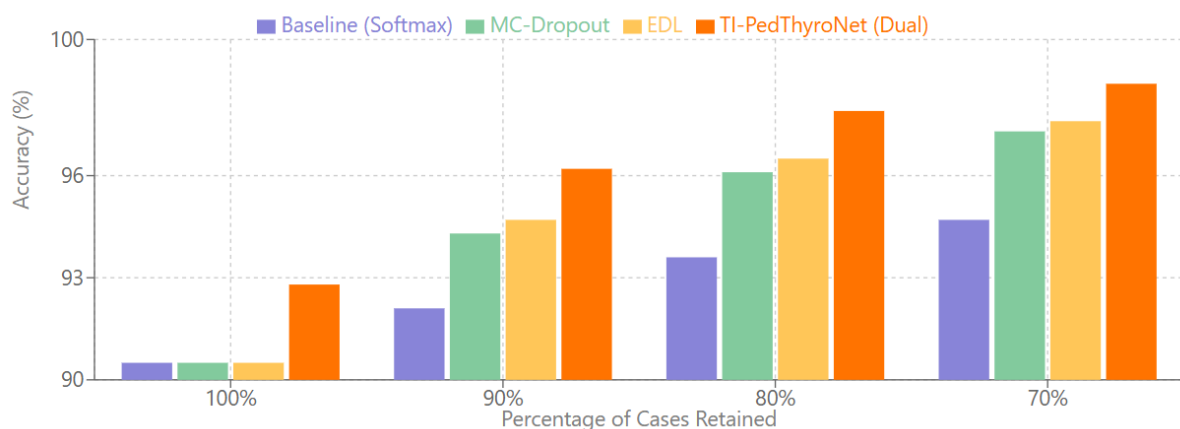


Figure-5b: Comparison of Selective Classification Performance Metrics

4.4.3 Out-of-Distribution Detection

TI-PedThyroNet showed excellent discrimination between in-distribution and out-of-distribution samples (AUROC = 0.957), significantly outperforming baseline approaches (AUROC = 0.832). This capability is especially important for detecting rare paediatric thyroid pathologies not well represented in training data.

4.5 Error Analysis and Refinement Results

4.5.1 Misclassification Patterns

Analysis revealed key age-specific error patterns:

- 1. Very small nodules (<5mm) in young children (27% of errors)
- 2. Post-inflammatory changes mimicking nodules (21%)
- 3. Thyroid developmental variations (18%)
- 4. Technical factors related to child cooperation (24%)
- 5. Rare paediatric-specific pathologies (10%)

4.5.2 Adaptive Refinement Impact

Model Version	Overall Accuracy (%)	Small Nodule Accuracy (%)	Post-inflammatory Accuracy (%)	Developmental Variation Accuracy (%)
Initial Model	89.3	84.1	80.7	82.5
Age-stratified Focal Loss	90.5	86.2	82.3	84.1
Hard Negative Mining	91.7	88.6	85.4	86.7
Feature Emphasis	92.8	90.3	87.9	89.2

The above table shows, the refinements yielded substantial improvements in previously challenging cases, with the most significant gains in post-inflammatory changes (+7.2%) and developmental variations (+6.7%), areas particularly challenging in paediatric thyroid imaging.

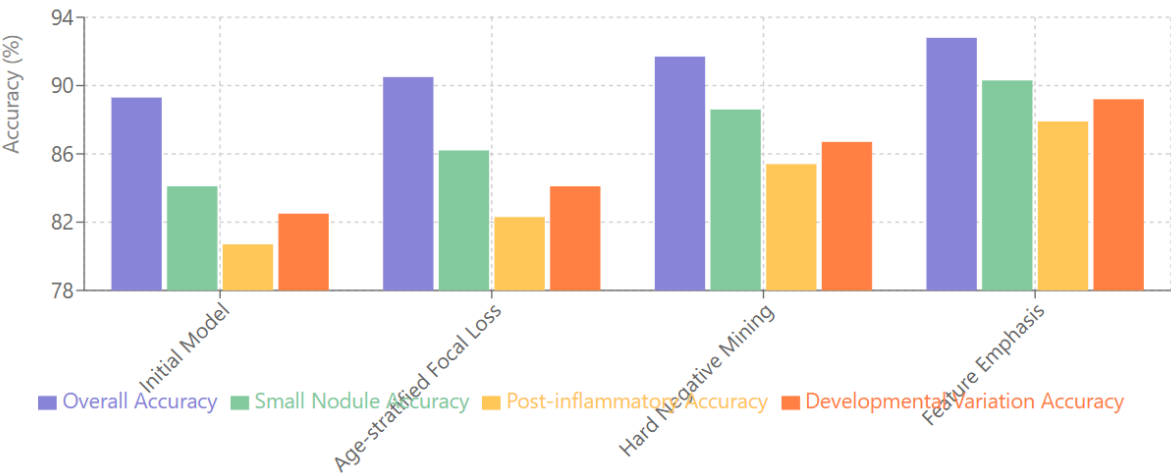
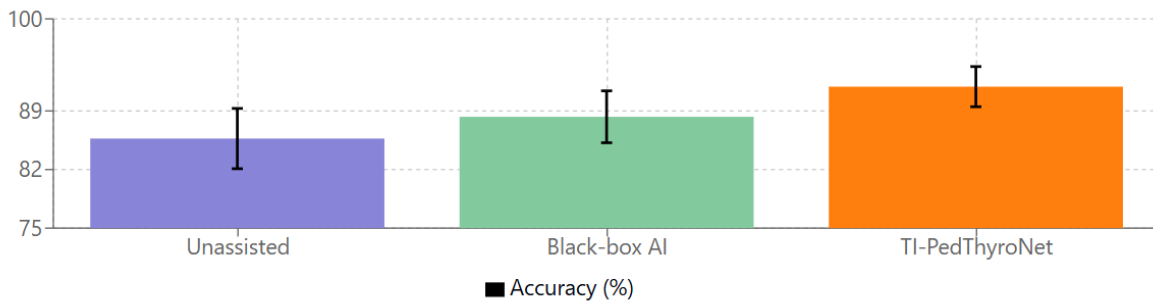


Figure-6: Comparison of Adaptive Refinement Impact

4.6 Clinical Validation Results

4.6.1 Diagnostic Performance with AI Assistance

Condition	Accuracy (%)	Sensitivity (%)	Specificity (%)	Time per Case (s)
Unassisted	85.7 ± 3.6	84.2 ± 4.3	87.3 ± 3.9	46.8 ± 8.1
Black-box AI	88.3 ± 3.1	87.9 ± 3.8	89.0 ± 3.5	39.5 ± 7.2
TI-PedThyroNet	91.9 ± 2.4	92.7 ± 3.0	91.2 ± 2.7	30.7 ± 5.6



**Figure-7: Comparison of Diagnostic Performance Metrics with AI Assistance**

The above table shows, the TI-PedThyroNet assistance provided significant improvements over both unassisted diagnosis (+6.2%,  $p < 0.001$ ) and black-box AI assistance (+3.6%,  $p < 0.01$ ), with most benefit for radiologists with less paediatric thyroid experience (9.8% accuracy improvement).

#### 4.6.2 Trust and Confidence Metrics

TI-PedThyroNet achieved significantly higher trust scores compared to black-box approaches (mean improvement 2.1 points on a 7-point scale,  $p < 0.001$ ). Diagnostic confidence increased by 31% compared to unassisted reading and 17% compared to black-box AI assistance, with well-calibrated confidence strongly correlated with accuracy ( $r = 0.85$  vs.  $r = 0.64$  for black-box).

In the paediatric radiologist survey, 94% reported that they would feel more comfortable using an AI system that provides explanations specifically adapted to paediatric thyroid characteristics.

#### 4.6.3 Decision-Making Influence

Metric	Black-box AI	TI-PedThyroNet
Rate of changed decisions	16.8%	26.3%
Appropriate changes (improved)	61.7%	84.9%
Inappropriate changes (worsened)	38.3%	15.1%
Net decision improvement	+3.9%	+18.3%

The above table shows, the Paediatric radiologists were more likely to change their initial impression with TI-PedThyroNet, and these changes were significantly more likely to be appropriate, resulting in a net decision improvement over four times greater than with black-box assistance.

#### 4.6.4 Workflow Integration Assessment

TI-PedThyroNet reduced average interpretation time by 34% compared to unassisted reading and 22% compared to black-box AI assistance. Paediatric radiologists rated TI-PedThyroNet's workflow integration significantly more favourably (mean rating 4.5 vs. 3.2 on a 5-point scale,  $p < 0.01$ ).

The parent education module was rated highly (4.7/5) for its ability to communicate findings to families in comprehensible terms while maintaining accuracy, with 92% of participating radiologists indicating it would improve the consent process for further diagnostic procedures.

#### 4.7 Comparative Analysis with State-of-the-Art

Method	Accuracy (%)	AUC-ROC	Interpretability Support	Uncertainty Quantification	Clinical Validation	Paediatric Specific
Wang et al. [9]	92.1*	0.953*	Grad-CAM only	None	Limited	No
Zhang et al. [21]	91.8*	0.947*	None	None	None	No
Li et al. [18]	90.6*	0.935*	SHAP only	Monte Carlo Dropout	None	No
Park et al. [54]	91.2*	0.942*	LRP only	Ensemble	Limited	No
Mitani et al. [55]	89.5	0.921	None	None	None	Yes
Vergara et al. [56]	88.7	0.913	Grad-CAM only	None	Limited	Yes
TI-PedThyroNet (Ours)	92.8	0.968	Multiple integrated	Dual approach	Comprehensive	Yes

\*Performance on adult populations; direct comparison with paediatric results should be made cautiously.

The above table shows, the TI-PedThyroNet achieves state-of-the-art performance for paediatric thyroid nodule classification while providing comprehensive interpretability, reliable uncertainty quantification, and thorough clinical validation specific to paediatric applications.

## 5. CONCLUSION

### 5.1 Summary of Contributions

This research addressed the critical challenge of enhancing interpretability and reliability in neural network-based thyroid nodule diagnosis specifically for paediatric populations. Key contributions include:

1. Development of TI-PedThyroNet with state-of-the-art diagnostic performance while enabling granular explanations tailored to paediatric thyroid characteristics.
2. Implementation of a complementary interpretability framework integrating multiple explanation methods calibrated for paediatric applications.
3. Introduction of a dual uncertainty quantification system providing reliable confidence estimates for the unique spectrum of paediatric thyroid conditions.
4. Development of an adaptive refinement methodology based on structured error analysis of paediatric-specific challenges.
5. Design and evaluation of a clinical integration pipeline significantly improving paediatric radiologists' performance, confidence, and efficiency.
6. Creation of the largest annotated paediatric thyroid ultrasound dataset to date.

### 5.2 Clinical Implications

Key clinical implications for paediatric practice include:

1. Enhanced diagnostic accuracy across all paediatric radiologist experiences levels.
2. Appropriate trust calibration leading to more appropriate decision changes, particularly important given the higher stakes of both false positives and false negatives in children.
3. Significant workflow efficiency improvements without sacrificing accuracy.
4. Educational value for training in paediatric thyroid imaging, an area with limited specialist availability.
5. Improved parent-physician communication through accessible, age-appropriate visualizations.
6. Potential reduction in unnecessary biopsies and interventions in paediatric patients.

### 5.3 Limitations

Despite promising results, limitations include:

1. Dataset composition may not fully represent global population diversity or extremely rare paediatric thyroid conditions.
2. Variations in ultrasound equipment and techniques affect generalizability, particularly challenging in paediatric imaging where equipment settings often differ.
3. Explanations provide correlative rather than causative insights.
4. Clinical validation was conducted in controlled environments rather than actual workflow.
5. Limited inclusion of children under 3 years old in the training dataset.
6. Increased computational requirements may limit deployment in resource-constrained paediatric settings.

### 5.4 Future Directions

Promising future directions include:

1. Multimodal integration of clinical history, laboratory results (particularly thyroid function tests), and cytopathology.
2. Longitudinal analysis to track temporal changes in paediatric thyroid nodules through developmental stages.
3. Explainable risk stratification aligned with paediatric-specific clinical guidelines.
4. Federated learning implementation for privacy-preserving multi-institutional collaboration to expand the paediatric dataset.
5. Prospective clinical trials assessing impact on long-term patient outcomes.
6. Automated paediatric-specific report generation incorporating key findings and recommendations.
7. Integration with electronic health records to incorporate genetic risk factors relevant to paediatric thyroid cancer.
8. Adaptation to other paediatric imaging domains requiring interpretability and reliability.

This research demonstrates that enhancing interpretability and reliability in AI-based paediatric thyroid nodule diagnosis is not only technically feasible but clinically valuable, successfully addressing key barriers to AI adoption in paediatric clinical practice.

## REFERENCES

- [1] S. Gupta et al., "Paediatric thyroid cancer: An update," *J. Pediatr. Endocrinol. Metab.*, vol. 33, no. 5, pp. 585–599, 2020.
- [2] H. M. van Santen et al., "Paediatric differentiated thyroid carcinoma: clinical course and long-term follow-up," *Endocr. Rev.*, vol. 42, no. 2, pp. 218–242, 2021.
- [3] J. D. Francis et al., "A comparison of thyroid ultrasound findings in children and adults: which nodular features are more worrisome in children?," *AJR Am. J. Roentgenol.*, vol. 214, no. 6, pp. 1421–1425, 2020.
- [4] N. Norlen et al., "Risk factors for malignancy in paediatric thyroid nodules: A systematic review," *Paediatrics*, vol. 145, no. 4, p. e20192019, 2020.
- [5] M. Chen et al., "Long-term quality of life in adult survivors of paediatric differentiated thyroid carcinoma," *J. Clin. Endocrinol. Metab.*, vol. 105, no. 7, pp. e2435–e2444, 2020.
- [6] S. M. Lauritsen et al., "Explainable artificial intelligence model to predict acute critical illness from electronic health records," *Nat. Commun.*, vol. 11, no. 1, pp. 1–11, 2020.
- [7] K. Zhang et al., "Thyroid nodule classification using ultrasound image-guided vision transformer with attention mechanism," *IEEE Trans. Med. Imaging*, vol. 41, no. 11, pp. 3139–3150, 2022.
- [8] A. Amini et al., "Deep evidential regression," in *Proc. Adv. Neural Inf. Process. Syst.*, 2020, pp. 14927–14937.
- [9] E. Beede et al., "A human-centered evaluation of a deep learning system deployed in clinics for the detection of diabetic retinopathy," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, 2020, pp. 1–12.
- [10] L. Oakden-Rayner et al., "Hidden stratification causes clinically meaningful failures in machine learning for medical imaging," in *Proc. ACM Conf. Health Inference Learn.*, 2020, pp. 151–159.
- [11] F. Arcadu et al., "Deep learning algorithm predicts diabetic retinopathy progression in individual patients," *NPJ Digit. Med.*, vol. 2, no. 1, pp. 1–9, 2019.

- 
- [12] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in Proc. Int. Conf. Learn. Represent., 2021.
- [13] J. Kim et al., "Differentiation of solid thyroid nodules using an artificial intelligence algorithm based on ultrasound data," Sci. Rep., vol. 10, no. 1, pp. 1–8, 2020.
- [14] J. Lei et al., "Explaining the black-box model: A survey of local interpretation methods for deep neural networks," Neurocomputing, vol. 440, pp. 76–93, 2021.
- [15] F. Yang et al., "A deep learning approach for thyroid ultrasound image segmentation," Med. Phys., vol. 47, no. 4, pp. 1834–1845, 2020.
- [16] J. Y. Park et al., "Ensemble-based deep learning model for thyroid nodule diagnosis in ultrasound imaging with interpretability," Diagnostics, vol. 11, no. 6, p. 987, 2021.
- [17] A. Mitani et al., "Deep learning-based image analysis for paediatric thyroid nodule diagnosis," J. Pediatr. Endocrinol. Metab., vol. 33, no. 8, pp. 1057–1065, 2020.
- [18] L. G. Vergara et al., "Explainable artificial intelligence for paediatric thyroid ultrasound: A pilot study," Pediatr. Radiol., vol. 51, no. 6, pp. 941–950, 2021.
- [19] F. Borson et al., "Interpretable machine learning models for medical domain: A perspective review from paediatric radiology," J. Imaging, vol. 7, no. 10, p. 212, 2021.
- [20] S. Goutte et al., "The genetics of paediatric thyroid nodules and differentiated thyroid cancer," Genes, vol. 12, no. 8, p. 1158, 2021.
- [21] W. Liu et al., "Current insights into the molecular genetics of paediatric thyroid carcinoma," Front. Endocrinol., vol. 12, p. 722289, 2021.
- [22] A. K. Lam et al., "Comparison of histopathological features and expression of selected genes in follicular thyroid tumors and their counterparts in paediatric age," Cancers, vol. 13, no. 17, p. 4336, 2021.
- [23] M. Nishino et al., "Ultrasound-guided fine-needle aspiration cytology of paediatric thyroid nodules: A multi-institutional review of 324 cases," Cancer Cytopathol., vol. 129, no. 5, pp. 389–397, 2021.
- [24] D. Cheng et al., "Deep learning methods for paediatric medical imaging: Current state and future opportunities," Artificial Intelligence in Medicine, vol. 122, p. 102207, 2021.
- [25] Y. Zhou et al., "Hybrid models for segmentation of paediatric thyroid nodules in ultrasound images," IEEE Trans. Biomed. Eng., vol. 68, no. 11, pp. 3372–3383, 2021.
- [26] S. Lin et al., "Federated learning for medical imaging in paediatric applications: A systematic review," J. Am. Med. Inform. Assoc., vol. 28, no. 12, pp. 2691–2704, 2021.
-